

# Text as Data

Justin Grimmer

Associate Professor  
Department of Political Science  
Stanford University

October 30th, 2014

# Structured Topic Models

## 1) Task:

- Examine how document attention, topic content varies↪ over time, across authors, or with **general set of covariates**

## 2) Objective Function

$$f(\mathbf{X}, \pi, \Theta, \alpha, \mathbf{W})$$

Where:

- $\mathbf{W}$  condition on information in document↪ other potential modifications to objective function. **Meta-data**.
- $f(\mathbf{X}, \pi, \Theta, \alpha, \mathbf{W})$  may encode additional information↪ layers of clustering, layers of topics, etc

## 3) Optimization

- EM, Variational Approximation, Gibbs Sampling, ...

## 4) Validation↪ many of the same methods from clustering

- Semantic, Convergent, Discriminant, Predictive, Hypothesis validity
- **How do we avoid the electric machine critique?**

# LDA Revisited

$$\boldsymbol{\theta}_k \sim \text{Dirichlet}(\mathbf{1})$$

$$\boldsymbol{\pi}_i | \boldsymbol{\alpha} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\tau_{im} | \boldsymbol{\pi}_i \sim \text{Multinomial}(1, \boldsymbol{\pi}_i)$$

$$x_{im} | \boldsymbol{\theta}_k, \tau_{imk} = 1 \sim \text{Multinomial}(1, \boldsymbol{\theta}_k)$$

# LDA Revisited

**Unigram Model<sub>k</sub>**  $\sim$  Dirichlet(**1**)

**Doc. Prop<sub>i</sub>**  $\sim$  Dirichlet(**Pop. Proportion**)

**Word Topic<sub>im</sub>**  $\sim$  Multinomial(1, **Doc. Prop<sub>i</sub>**)

**Word<sub>im</sub>**  $\sim$  Multinomial(1, **Unigram Model<sub>k</sub>**)

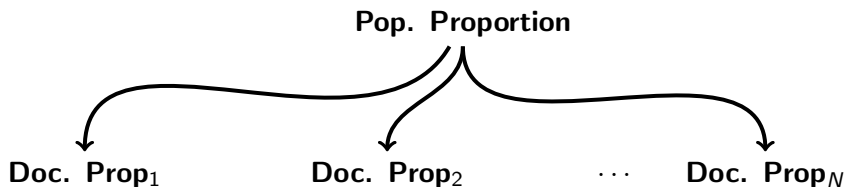
# A General Hierarchical Structure

LDA:

**Pop. Proportion**

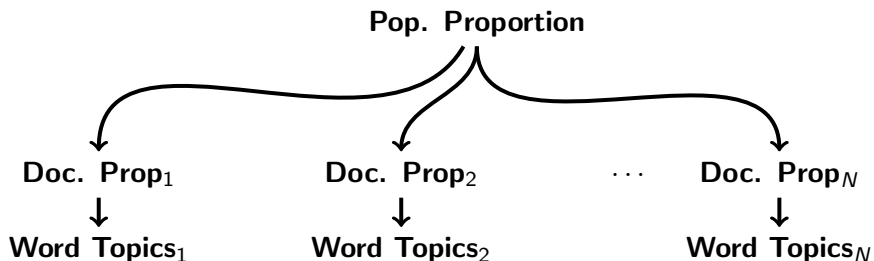
# A General Hierarchical Structure

LDA:



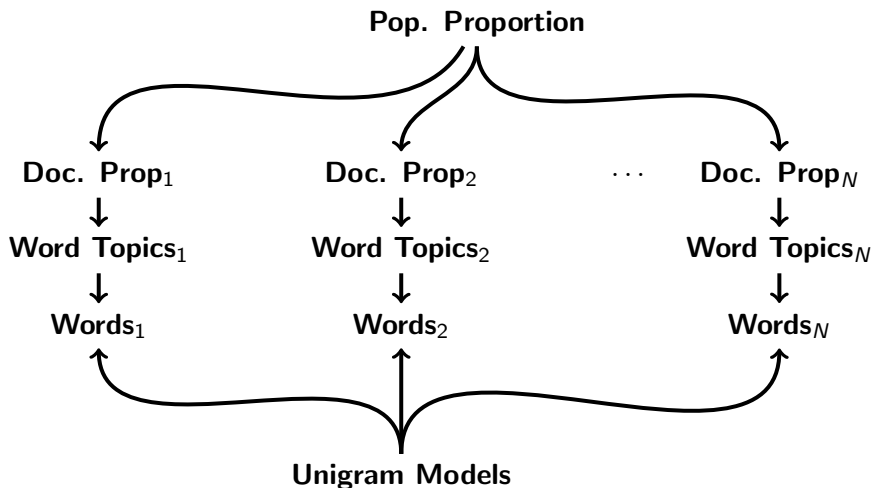
# A General Hierarchical Structure

LDA:



# A General Hierarchical Structure

LDA:





# A General Hierarchical Structure

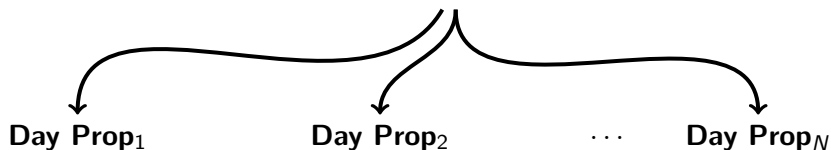
Dynamic Topic Model (Quinn et al 2010)

**Dynamic Prior Across Days**

# A General Hierarchical Structure

Dynamic Topic Model (Quinn et al 2010)

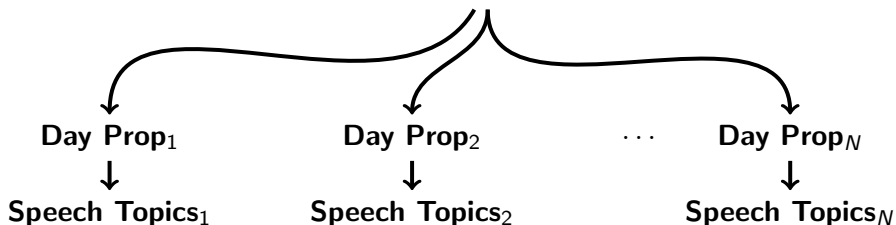
**Dynamic Prior Across Days**



# A General Hierarchical Structure

Dynamic Topic Model (Quinn et al 2010)

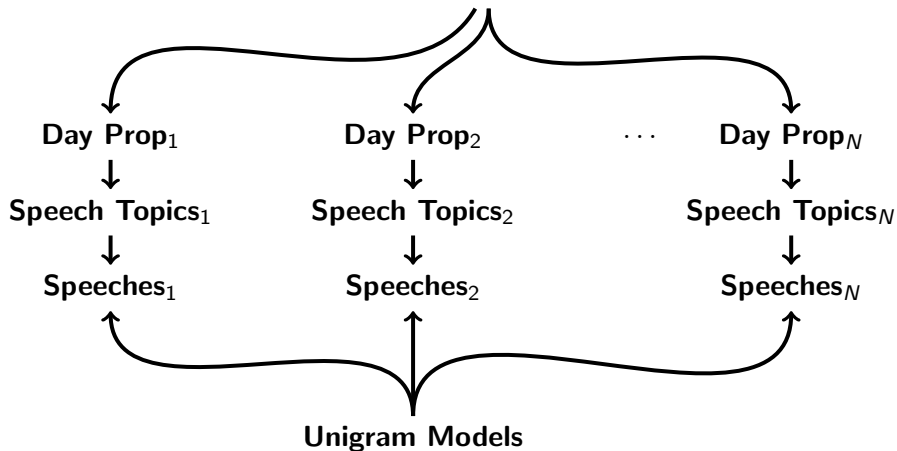
**Dynamic Prior Across Days**



# A General Hierarchical Structure

Dynamic Topic Model (Quinn et al 2010)

**Dynamic Prior Across Days**



# A General Hierarchical Structure

Expressed Agenda Model (Grimmer 2010)

**Average Attention Across Authors**

# A General Hierarchical Structure

Expressed Agenda Model (Grimmer 2010)

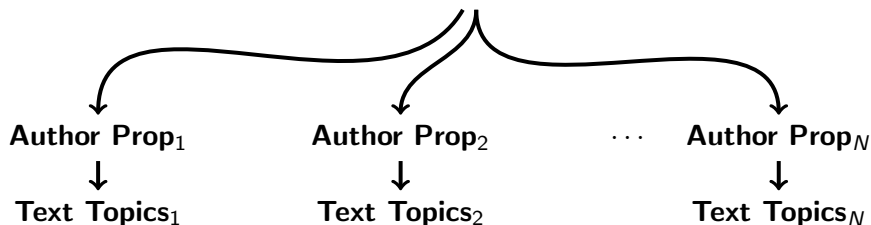
**Average Attention Across Authors**



# A General Hierarchical Structure

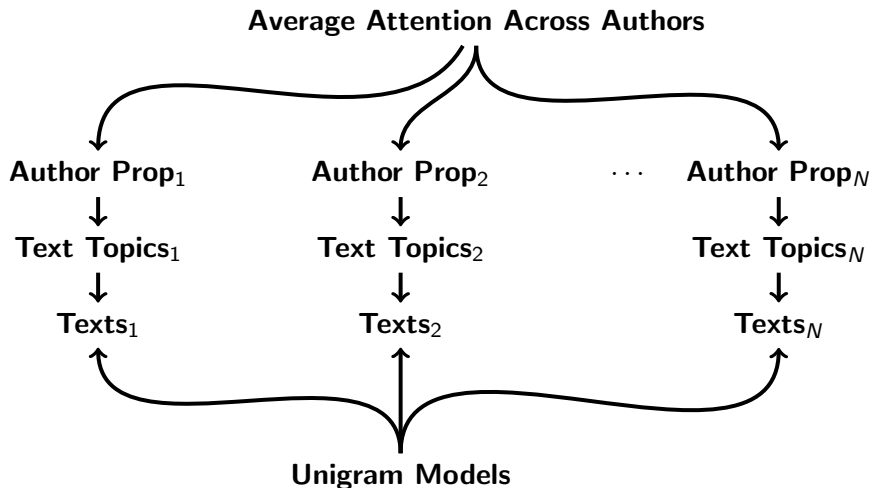
Expressed Agenda Model (Grimmer 2010)

**Average Attention Across Authors**



# A General Hierarchical Structure

Expressed Agenda Model (Grimmer 2010)





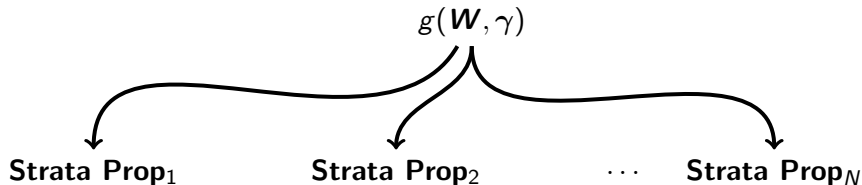
# A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)

$$g(\mathbf{W}, \gamma)$$

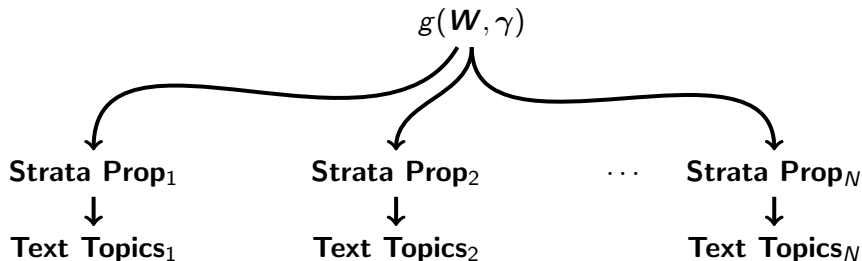
# A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



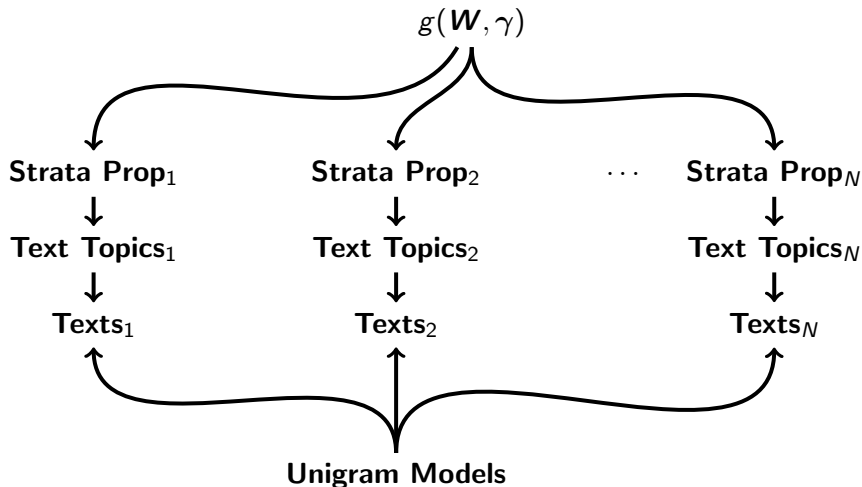
# A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



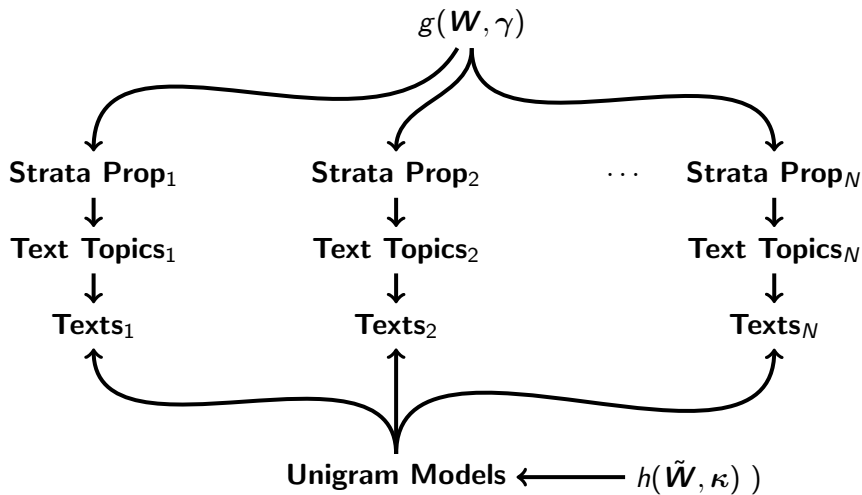
# A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



# A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



# A General Hierarchical Structure

Conditioning on Unknown Covariates  $\rightsquigarrow$  levels of mixtures at proportions  
(Grimmer 2013; Wallach 2008)

## Mixture of Top. Attn. Models

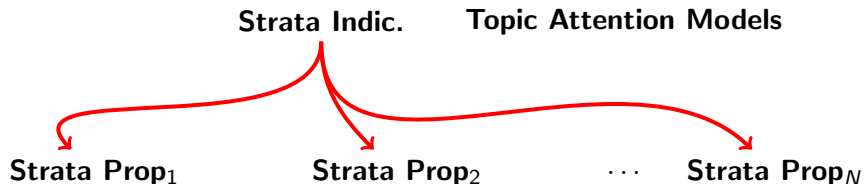
# A General Hierarchical Structure

Conditioning on Unknown Covariates  $\rightsquigarrow$  levels of mixtures at proportions  
(Grimmer 2013; Wallach 2008)

**Strata Indic.**      **Topic Attention Models**

# A General Hierarchical Structure

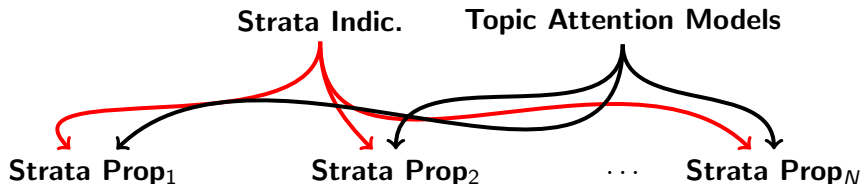
Conditioning on Unknown Covariates  $\rightsquigarrow$  levels of mixtures at proportions  
(Grimmer 2013; Wallach 2008)





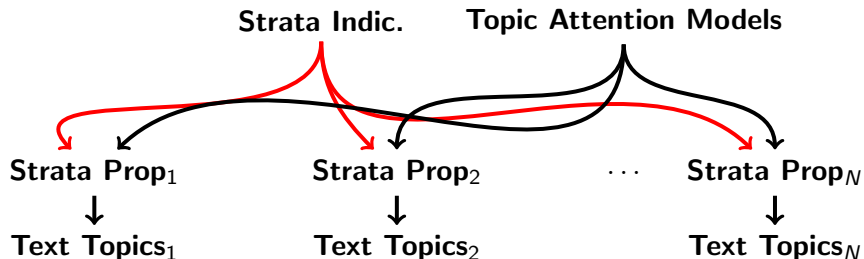
# A General Hierarchical Structure

Conditioning on Unknown Covariates  $\rightsquigarrow$  levels of mixtures at proportions  
(Grimmer 2013; Wallach 2008)



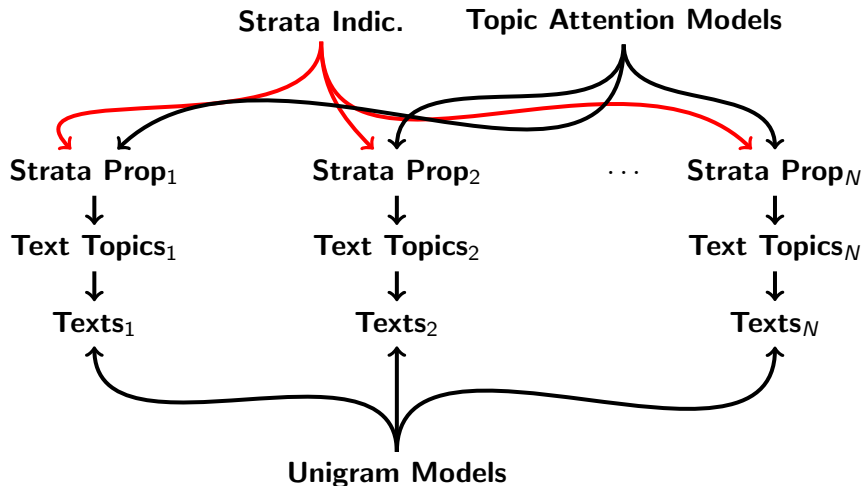
# A General Hierarchical Structure

Conditioning on Unknown Covariates  $\rightsquigarrow$  levels of mixtures at proportions  
(Grimmer 2013; Wallach 2008)



# A General Hierarchical Structure

Conditioning on Unknown Covariates  $\rightsquigarrow$  levels of mixtures at proportions  
(Grimmer 2013; Wallach 2008)



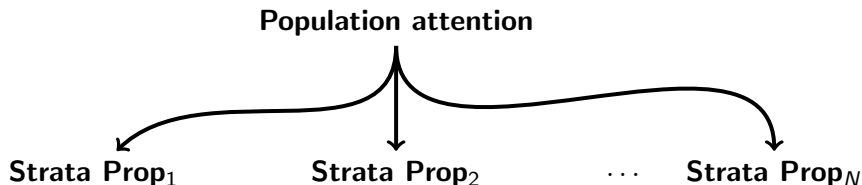
# A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics  $\rightsquigarrow$  hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2014)

**Population attention**

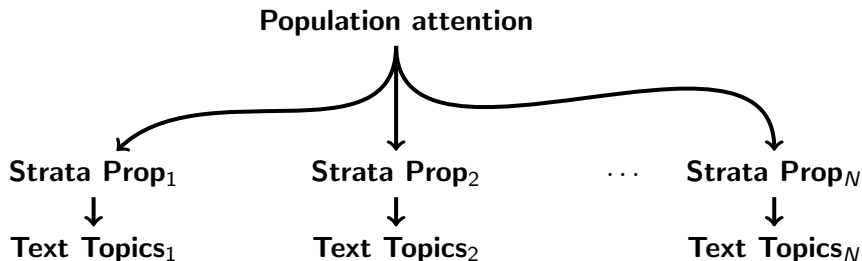
# A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics  $\rightsquigarrow$  hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2014)



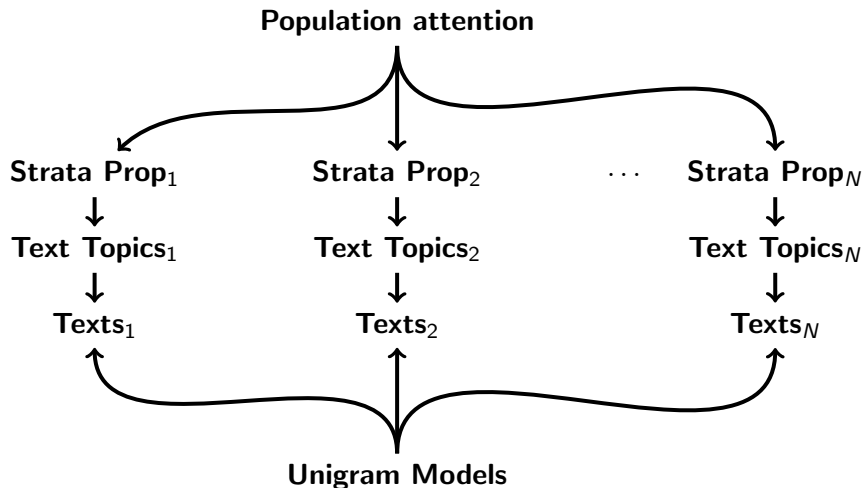
# A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics  $\rightsquigarrow$  hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2014)



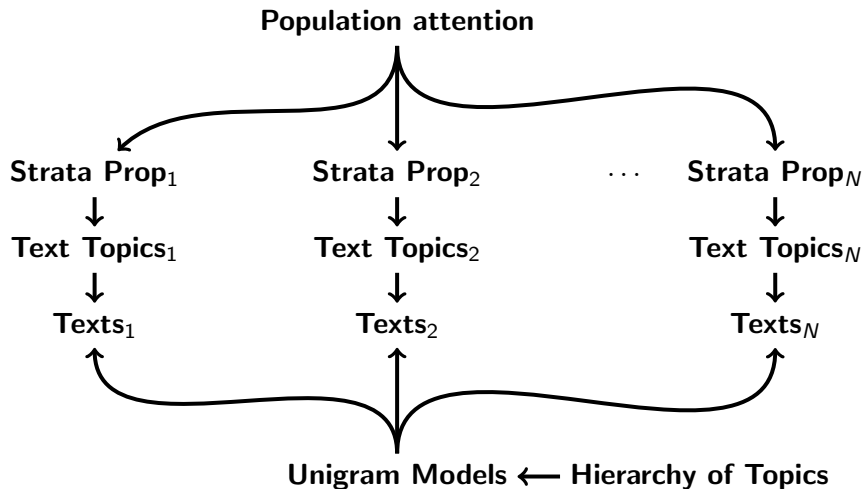
# A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics  $\rightsquigarrow$  hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2014)



# A General Hierarchical Structure

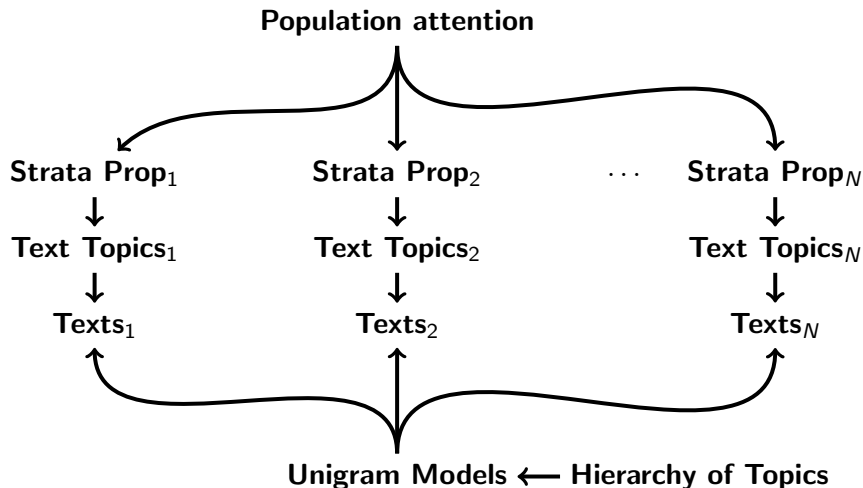
Conditioning on Unknown Covariates for Topics  $\rightsquigarrow$  hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2014)





# A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics  $\rightsquigarrow$  hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2014)



# Why Encode Structure in Extensions of LDA?

# Why Encode Structure in Extensions of LDA?

- Substantive reasons

# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content

# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content
  - Avoids potential ad-hoc secondary analysis

# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content
  - Avoids potential ad-hoc secondary analysis
  - Clear data generating process

# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content
  - Avoids potential ad-hoc secondary analysis
  - Clear data generating process
- Statistical reasons

# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content
  - Avoids potential ad-hoc secondary analysis
  - Clear data generating process
- Statistical reasons
  - **Smoothing**  $\rightsquigarrow$  borrow information across groups intelligently



# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content
  - Avoids potential ad-hoc secondary analysis
  - Clear data generating process
- Statistical reasons
  - **Smoothing**  $\rightsquigarrow$  borrow information across groups intelligently
  - **Uncertainty**  $\rightsquigarrow$  potential for better uncertainty estimates

# Why Encode Structure in Extensions of LDA?

- Substantive reasons
  - Additional structure corresponds to substantively interesting content
  - Avoids potential ad-hoc secondary analysis
  - Clear data generating process
- Statistical reasons
  - **Smoothing**  $\rightsquigarrow$  borrow information across groups intelligently
  - **Uncertainty**  $\rightsquigarrow$  potential for better uncertainty estimates
  - **Improved topics**  $\rightsquigarrow$  small word conditions, structure could help

# Plan for the Class

- 1) Discuss model with unknown covariates for strata proportions  $\rightsquigarrow$  presentational style
- 2) Discuss model with hierarchy of topics  $\rightsquigarrow$  mirrors genre

# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

- Characterize representation provided to constituents



# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

- Characterize representation provided to constituents
- Divide attention over a set of topics

# Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

- Characterize representation provided to constituents
- Divide attention over a set of topics
- Given attention to topics, write press releases

# Presentational Styles $\rightsquigarrow$ Objective Function

- $\pi_{itk} \equiv$  Attention senator  $i$  allocates to issue  $k$  in year  $t$
- $\pi_{itk} \equiv$  Probability press release is about issue  $k$
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

# Presentational Styles $\rightsquigarrow$ Objective Function

- $\pi_{itk} \equiv$  Attention senator  $i$  allocates to issue  $k$  in year  $t$
- $\pi_{itk} \equiv$  Probability press release is about issue  $k$
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release  $j$  from senator  $i$  in year  $t$ )

# Presentational Styles $\rightsquigarrow$ Objective Function

- $\pi_{itk} \equiv$  Attention senator  $i$  allocates to issue  $k$  in year  $t$
- $\pi_{itk} \equiv$  Probability press release is about issue  $k$
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release  $j$  from senator  $i$  in year  $t$ )

- **Assume**: Each press release  $j$  assigned to one topic.
- Let  $\tau_{ijt}$  indicate press release  $j$ 's topic.

# Presentational Styles $\rightsquigarrow$ Objective Function

- $\pi_{itk} \equiv$  Attention senator  $i$  allocates to issue  $k$  in year  $t$
- $\pi_{itk} \equiv$  Probability press release is about issue  $k$
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release  $j$  from senator  $i$  in year  $t$ )

- **Assume**: Each press release  $j$  assigned to one topic.
- Let  $\tau_{ijt}$  indicate press release  $j$ 's topic.

$$\tau_{ijt} \sim \text{Multinomial}(1, \boldsymbol{\pi}_{it})$$

# Presentational Styles $\rightsquigarrow$ Objective Function

- $\pi_{itk} \equiv$  Attention senator  $i$  allocates to issue  $k$  in year  $t$
- $\pi_{itk} \equiv$  Probability press release is about issue  $k$
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release  $j$  from senator  $i$  in year  $t$ )

- **Assume**: Each press release  $j$  assigned to one topic.
- Let  $\tau_{ijt}$  indicate press release  $j$ 's topic.

$$\tau_{ijt} \sim \text{Multinomial}(1, \boldsymbol{\pi}_{it})$$

- Conditional on topic, draw document's content.

# Presentational Styles $\rightsquigarrow$ Objective Function

- $\pi_{itk} \equiv$  Attention senator  $i$  allocates to issue  $k$  in year  $t$
- $\pi_{itk} \equiv$  Probability press release is about issue  $k$
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release  $j$  from senator  $i$  in year  $t$ )

- **Assume**: Each press release  $j$  assigned to one topic.
- Let  $\tau_{ijt}$  indicate press release  $j$ 's topic.

$$\tau_{ijt} \sim \text{Multinomial}(1, \boldsymbol{\pi}_{it})$$

- Conditional on topic, draw document's content.
- If  $\tau_{ijtk} = 1$  then

$$\mathbf{x}_{ijt} \sim \text{Multinomial}(n_{ijt}, \boldsymbol{\theta}_k).$$



# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

$$\sigma_{it} \sim \text{Multinomial}(1, \beta).$$

# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s)\end{aligned}$$

# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Other priors:

# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Other priors:

$$\theta_k \sim \text{Multinomial}(\lambda)$$

# Priors

Each  $\pi_{it}$  is a draw from one-of- $S$  styles  $\rightsquigarrow$  mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Other priors:

$$\begin{aligned}\theta_k &\sim \text{Multinomial}(\lambda) \\ \beta &\sim \text{Multinomial}(\mathbf{1})\end{aligned}$$

# Presentational Styles $\rightsquigarrow$ Objective Function



# Presentational Styles $\rightsquigarrow$ Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\boldsymbol{\lambda}) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

# Presentational Styles $\rightsquigarrow$ Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta)\end{aligned}$$

# Presentational Styles $\rightsquigarrow$ Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s)\end{aligned}$$

# Presentational Styles $\rightsquigarrow$ Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \tau_{ijt} | \pi_{it} &\sim \text{Multinomial}(1, \pi_{it})\end{aligned}$$

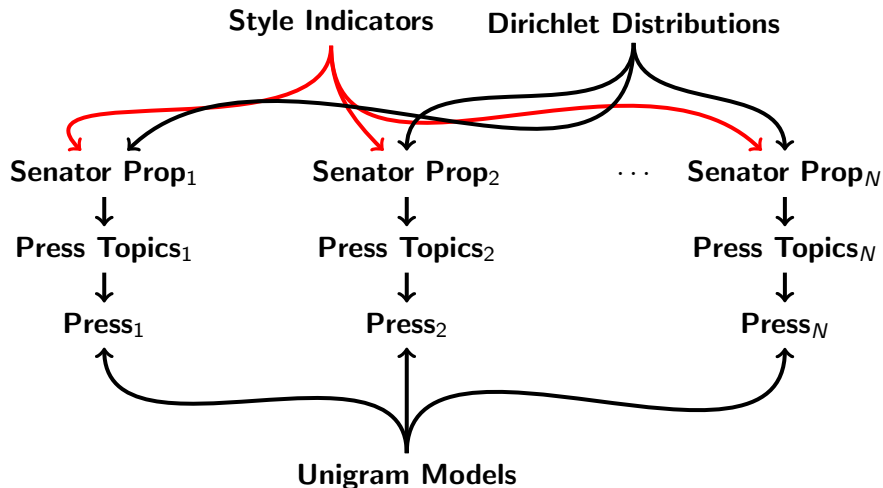
# Presentational Styles $\rightsquigarrow$ Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \tau_{ijt} | \pi_{it} &\sim \text{Multinomial}(1, \pi_{it}) \\ \mathbf{x}_{ijt} | \tau_{ijtk} = 1, \theta_k &\sim \text{Multinomial}(n_{ijt}, \theta_k)\end{aligned}$$

# Presentational Styles $\rightsquigarrow$ Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \tau_{ijt} | \pi_{it} &\sim \text{Multinomial}(1, \pi_{it}) \\ \mathbf{x}_{ijt} | \tau_{ijtk} = 1, \theta_k &\sim \text{Multinomial}(n_{ijt}, \theta_k)\end{aligned}$$

# Mixture of Styles, Mixture of Topics



## Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[ \beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[ \pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$



Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[ \beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[ \pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

## 1) Estimate with Variational Approximation

Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[ \beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[ \pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

- 1) Estimate with Variational Approximation
- 2) Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)

Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[ \beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[ \pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

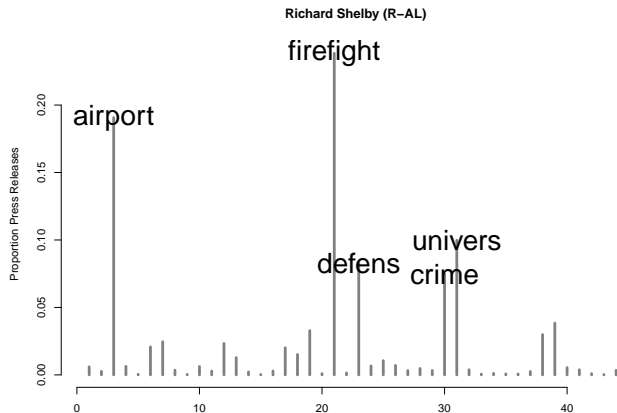
- 1) Estimate with Variational Approximation
- 2) Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)
  - Non-parametric model  $\rightsquigarrow$  statistical selection

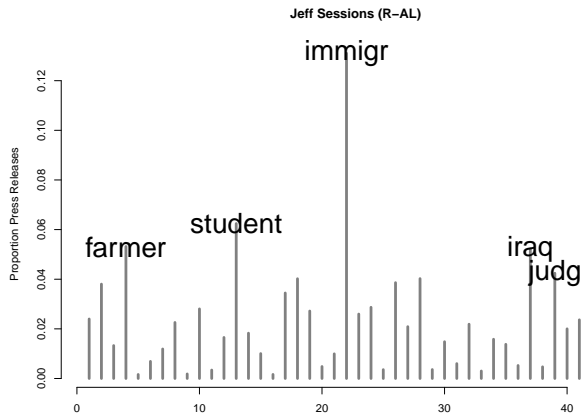
Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[ \beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[ \pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

- 1) Estimate with Variational Approximation
- 2) Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)
  - Non-parametric model  $\rightsquigarrow$  statistical selection
  - Experiments/Coding Exercises to assess





# Notions of validity: From Quinn, Monroe, et al (2010)

- **Semantic Validity:** All categories are coherent and meaningful
- **Convergent Construct Validity:** Measures concur with existing measures in critical details.
- **Discriminant Construct Validity:** Measures differ from existing measures in productive ways.
- **Predictive Measure:** Measures from the model corresponds to external events in expected ways.
- **Hypothesis Validity:** Measures generated from the model can be used to test substantive hypotheses.

To establish utility of new measures, demonstrate variety of **validations**

**None of these validations are performed using a canned statistic**

**All:** require substantive knowledge on areas (and what we expect!) [

# Home Style Measures, Semantic Validity

**Must:** Demonstrate to reader that topics are coherent and semantically meaningful

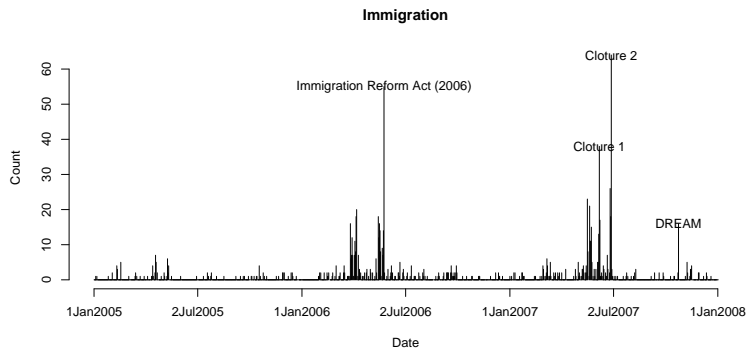
Description	Stems	%
Honorary	honor,prayer,rememb,fund,tribut	5.0
Transp. Grants	airport,transport,announc,urban,hud	4.8
Iraq	iraq,iraqi,troop,war,sectarian	4.7
DHS Policy	homeland,port,terrorist,dh,fema	4.1
Judicial Nom.	judg,court,suprem,nomin,nomine	3.8
Fire Dept. Grant	firefight,homeland,afgp,award,equip	3.7

How: **examples** in text are also useful.



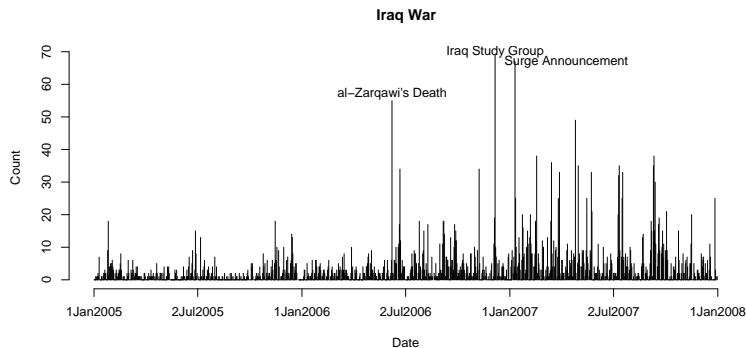
# Home Style Measures, Convergent Validity

## Over time variation



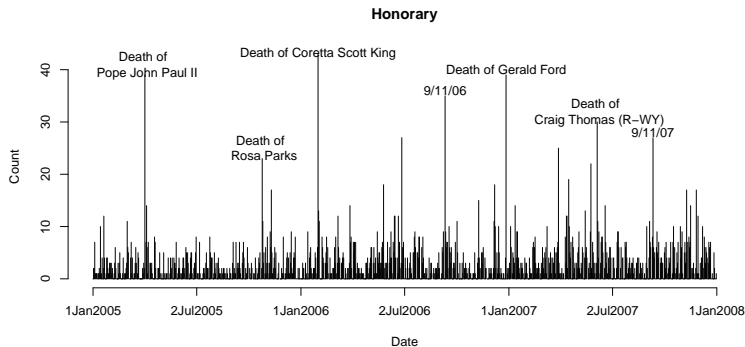
# Home Style Measures, Convergent Validity

## Over time variation



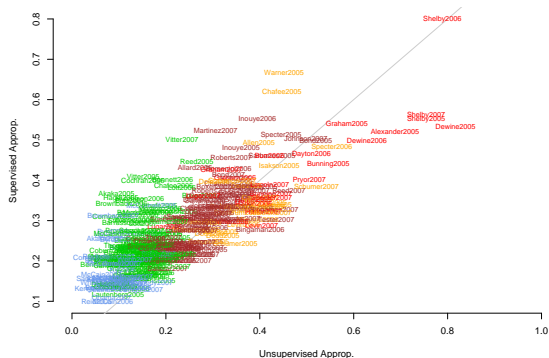
# Home Style Measures, Convergent Validity

## Over time variation

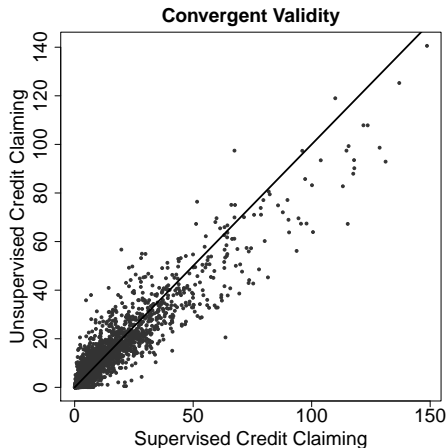


# Home Style Measures, Convergent Validity

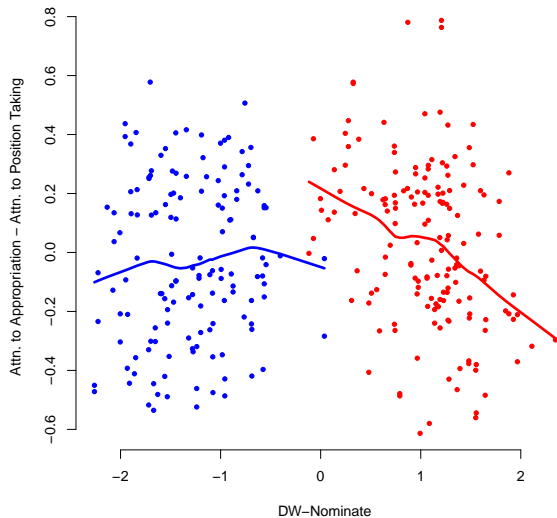
## Supervised/Unsupervised Convergence



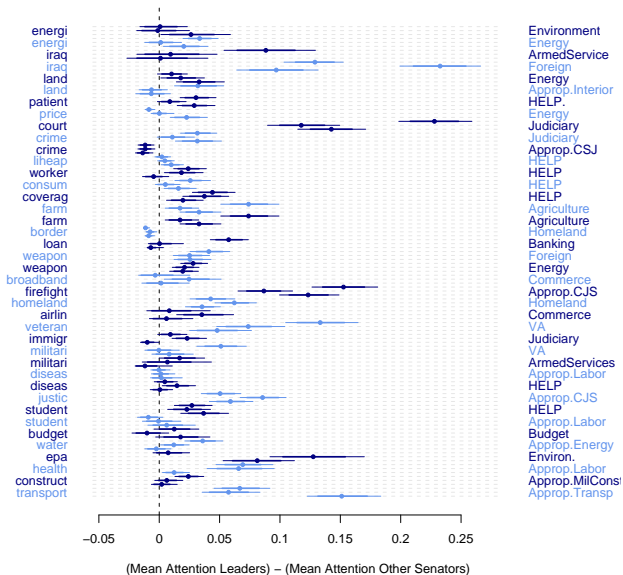
# Home Style Measures, Convergent Validity



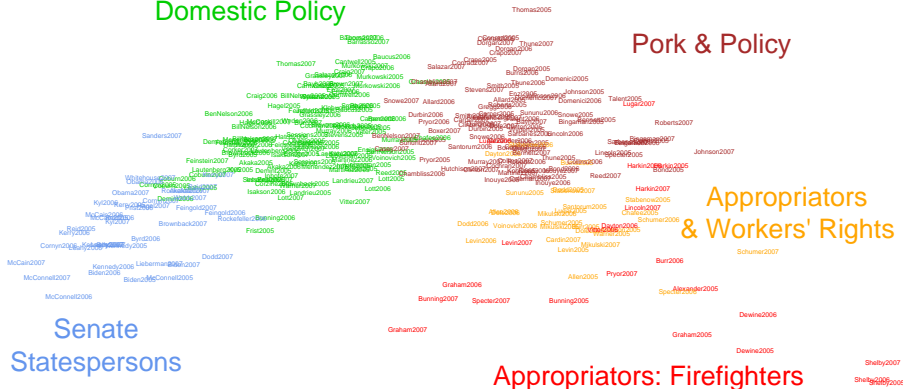
# Discriminant Construct Validity



# Predictive Validity



## Domestic Policy





# Hypothesis Validity

## Domestic Policy

## Pork & Policy

## Appropriators & Workers' Rights

## Appropriators: Firefighters

## Senate Statespersons

## Senate Statesperson

- Iraq War
- Intelligence
- Intl.  
Relations

# Hypothesis Validity

## Domestic Policy

## Pork & Policy

## Appropriators & Workers' Rights

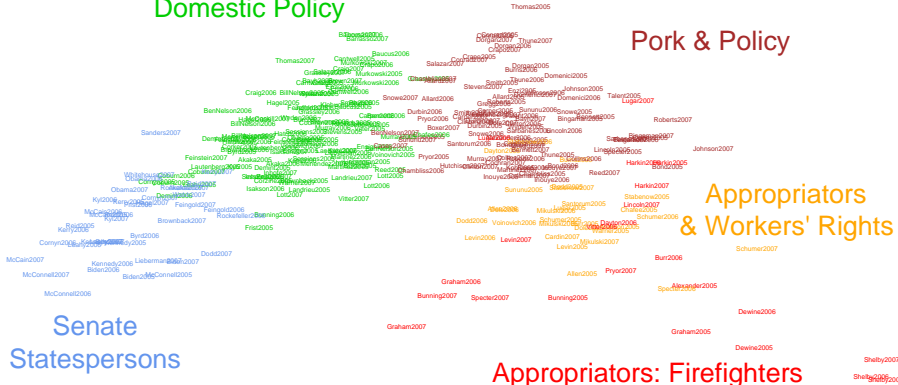
## Appropriators: Firefighters

## Senate Statespersons

## Domestic Policy

- Iraq War
- Intelligence
- Intl.  
Relations

- Environment
- Gas prices
- DHS
- Consumer



# Hypothesis Validity

## Domestic Policy

## Pork & Policy

## Appropriators & Workers' Rights

## Appropriators: Firefighters

## Senate Statespersons

## Domestic Policy

## Pork & Policy

- Iraq War
- Intelligence
- Intl. Relations

- Environment
- Gas prices
- DHS
- Consumer

- WRDA grants
- Farming
- Health Care
- Education

# Hypothesis Validity

## Domestic Policy

## Pork & Policy

## Appropriators & Workers' Rights

## Appropriators: Firefighters

## Senate Statespersons

## Senate Statesperson

- Iraq War
- Intelligence
- Intl.  
Relations

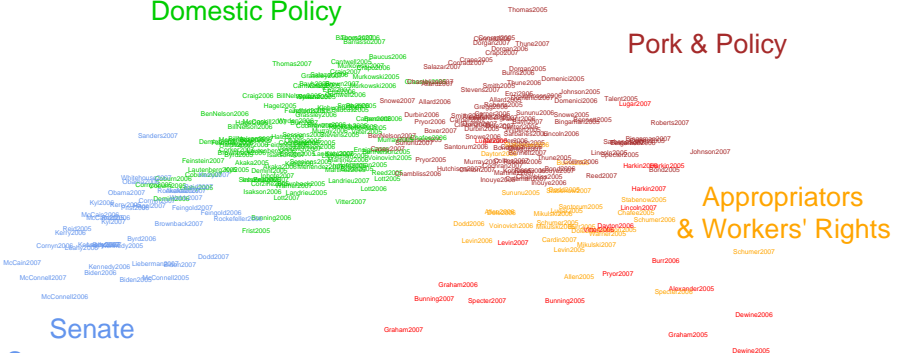
- ## Domestic Policy
- Environment
  - Gas prices
  - DHS
  - Consumer

## Pork & Policy

- WRDA  
grants
- Farming
- Health Care
- Education

## Appropriators

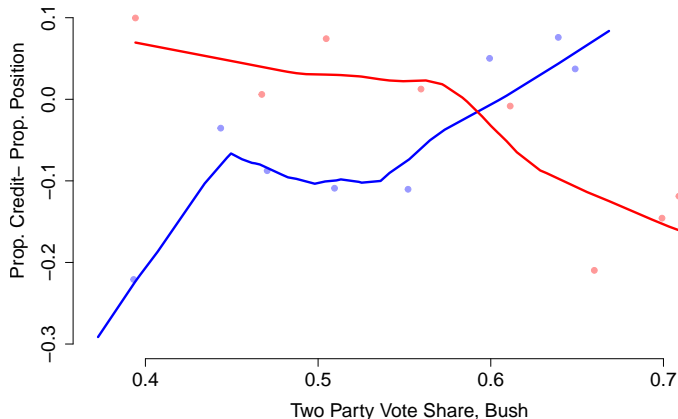
- Fire Grants
- Airport  
Grants
- University
- Money



# Hypothesis Validity

Why do senators adopt different styles?

## District Fit



# What are the right number of topics?

# What are the right number of topics?

- Number of topics  $\rightsquigarrow$  depends on task at hand

# What are the right number of topics?

- Number of topics  $\rightsquigarrow$  depends on task at hand
- Coarse  $\rightsquigarrow$  broad comparisons, lose distinctions



# What are the right number of topics?

- Number of topics  $\rightsquigarrow$  depends on task at hand
- Coarse  $\rightsquigarrow$  broad comparisons, lose distinctions
- Granular  $\rightsquigarrow$  specific insights, lose broader picture

# What are the right number of topics?

- Number of topics  $\rightsquigarrow$  depends on task at hand
- Coarse  $\rightsquigarrow$  broad comparisons, lose distinctions
- Granular  $\rightsquigarrow$  specific insights, lose broader picture
- **Hierarchy of topics**  $\rightsquigarrow$  Pachinko Allocation, Hierarchies of von-Mises Fisher Distributions

# What are the right number of topics?

- Number of topics  $\rightsquigarrow$  depends on task at hand
- Coarse  $\rightsquigarrow$  broad comparisons, lose distinctions
- Granular  $\rightsquigarrow$  specific insights, lose broader picture
- **Hierarchy of topics**  $\rightsquigarrow$  Pachinko Allocation, Hierarchies of von-Mises Fisher Distributions

Blaydes, Grimmer, and McQueen [In Progress]  $\rightsquigarrow$  estimate nested topics to explore the **Mirros for Princes**

# The Mirrors Genre

26 Christian mirrors

# The Mirrors Genre

26 Christian mirrors

- The Prince (1513 CE)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts



# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

## Work with translations

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations~> little evidence of selection

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations~> little evidence of selection

- Collect data on collection of 98 (51 Christian, 47 Islamic, some not translated)

# The Mirrors Genre

## 26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

## 21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations~> little evidence of selection

- Collect data on collection of 98 (51 Christian, 47 Islamic, some not translated)
- No difference on Year/Region

# Preprocessing Texts

47 books



# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (**allah**) and discard proper nouns
- Identified synonyms
  - **almighty, god**

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (**allah**) and discard proper nouns
- Identified synonyms
  - **almighty, god**
  - **monarch, prince, king, ruler**

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
  - almighty, god
  - monarch, prince, king, ruler
  - Lord  $\neq$  lord



# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (**allah**) and discard proper nouns
- Identified synonyms
  - **almighty**, **god**
  - **monarch**, **prince**, **king**, **ruler**
  - **Lord**  $\neq$  **lord**

Result: short segment  $j$  in book  $i$  is a count vector

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
  - almighty, god
  - monarch, prince, king, ruler
  - Lord  $\neq$  lord

Result: short segment  $j$  in book  $i$  is a count vector

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ij2124})$$

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
  - almighty, god
  - monarch, prince, king, ruler
  - Lord  $\neq$  lord

Result: short segment  $j$  in book  $i$  is a count vector

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ij2124})$$

We work with a normalized version of the documents,

# Preprocessing Texts

47 books  $\rightsquigarrow$  Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
  - almighty, god
  - monarch, prince, king, ruler
  - Lord  $\neq$  lord

Result: short segment  $j$  in book  $i$  is a count vector

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ij2124})$$

We work with a normalized version of the documents,

$$\mathbf{x}_{ij}^* = \frac{\mathbf{x}_{ij}}{\sqrt{\mathbf{x}_{ij}' \mathbf{x}_{ij}}}$$

# Measuring Themes in the Mirrors

Model built around two hierarchies:

# Measuring Themes in the Mirrors

Model built around two hierarchies:

- 1) Books  $\rightsquigarrow$  paragraphs (Blei, Ng, Jordan 2003; Wallach, 2008; Quinn et al 2010; Grimmer 2010; Roberts et al 2014)

# Measuring Themes in the Mirrors

Model built around two hierarchies:

- 1) Books  $\rightsquigarrow$  paragraphs (Blei, Ng, Jordan 2003; Wallach, 2008; Quinn et al 2010; Grimmer 2010; Roberts et al 2014)
- 2) Coarse topics  $\rightsquigarrow$  granular topics (Li and McCallum 2006; Gopal and Yang 2014)

# Measuring Themes in the Mirrors

Estimate **four** quantities of interest



# Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)

# Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)

# Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)
  - Each granular topic classified into one coarse topic

# Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)
  - Each granular topic classified into one coarse topic
- 3) Each book  $i$ 's **themes** <sub>$i$</sub>

$$\mathbf{themes}_i = (\text{theme}_{i1}, \text{theme}_{i2}, \dots, \text{theme}_{i60})$$

# Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)
  - Each granular topic classified into one coarse topic
- 3) Each book *i*'s **themes**;
- 4) Each short segment's granular (and coarse) topic

# A Hierarchy of Topics

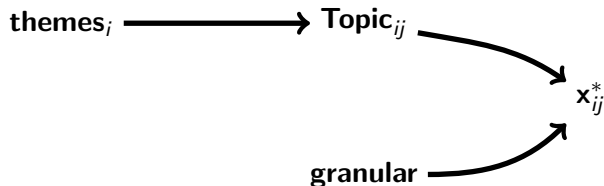
**themes;**

# A Hierarchy of Topics

**themes<sub>*i*</sub>  $\longrightarrow$  Topic<sub>*ij*</sub>**

$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

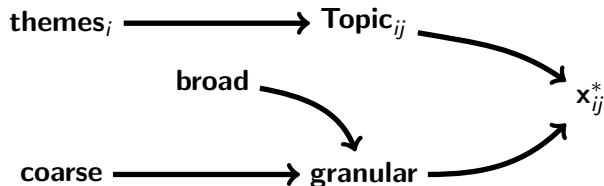
# A Hierarchy of Topics



$$\begin{aligned} \mathbf{Topic}_{ij} &\sim \text{Multinomial}(1, \mathbf{themes}_i) \\ \mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 &\sim \text{vMF}(\kappa, \mathbf{granular}_k) \end{aligned}$$



# A Hierarchy of Topics



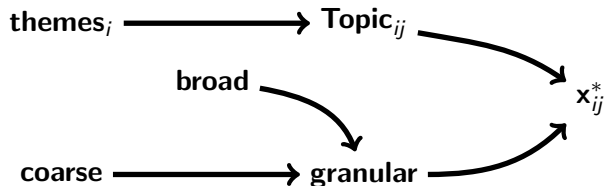
$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

$$\mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 \sim \text{vMF}(\kappa, \mathbf{granular}_k)$$

$$\mathbf{broad}_k \sim \text{Multinomial}(1, \mathbf{Broad Theme Prior})$$

$$\mathbf{granular}_k | \mathbf{broad}_{km} = 1 \sim \text{vMF}(\kappa, \mathbf{coarse}_m)$$

# A Hierarchy of Topics



$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

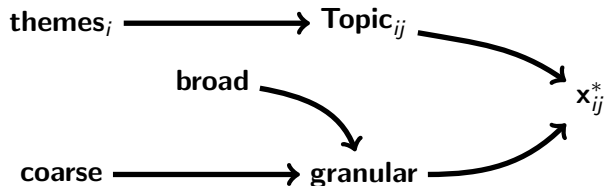
$$\mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 \sim \text{vMF}(\kappa, \mathbf{granular}_k)$$

$$\mathbf{broad}_k \sim \text{Multinomial}(1, \mathbf{Broad Theme Prior})$$

$$\mathbf{granular}_k | \mathbf{broad}_{km} = 1 \sim \text{vMF}(\kappa, \mathbf{coarse}_m)$$

Estimate model with Variational Approximation

# A Hierarchy of Topics



$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

$$\mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 \sim \text{vMF}(\kappa, \mathbf{granular}_k)$$

$$\mathbf{broad}_k \sim \text{Multinomial}(1, \mathbf{Broad Theme Prior})$$

$$\mathbf{granular}_k | \mathbf{broad}_{km} = 1 \sim \text{vMF}(\kappa, \mathbf{coarse}_m)$$

Estimate model with Variational Approximation

Model selection: automatic model fit, qualitative evaluation

# Interpreting Unsupervised Models

Two approaches to labeling output

# Interpreting Unsupervised Models

Two approaches to labeling output

- 1) **Computational**: identify discriminating words

# Interpreting Unsupervised Models

Two approaches to labeling output

- 1) **Computational**: identify discriminating words
- 2) **Manual**: Segments classified to coarse, granular topics. Read, discuss, and label

# Interpreting Unsupervised Models

Two approaches to labeling output

- 1) **Computational**: identify discriminating words
- 2) **Manual**: Segments classified to coarse, granular topics. Read, discuss, and label

Unsupervised models **structure** and **guide** our reading

# Art of Rulership

Practices and ideals of political rule



# Art of Rulership

Practices and ideals of political rule

king

# Art of Rulership

Practices and ideals of political rule

king, princ

# Art of Rulership

Practices and ideals of political rule

king, princ, citi

# Art of Rulership

Practices and ideals of political rule

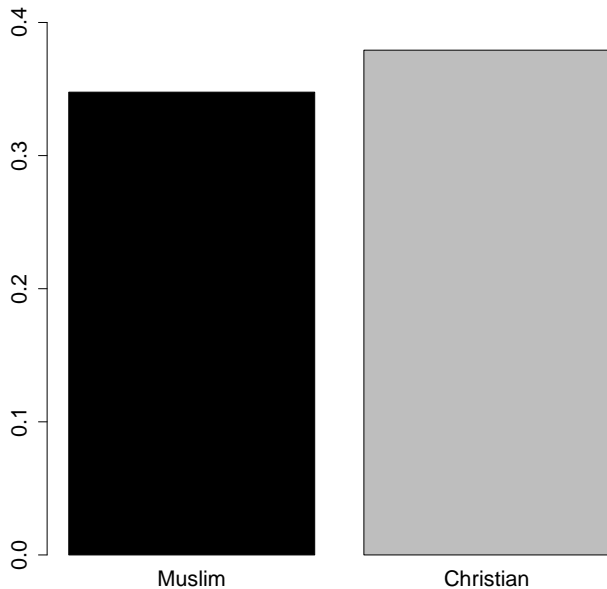
king, princ, citi, great, place, work, emperor, enemi, armi, letter

# Art of Rulership

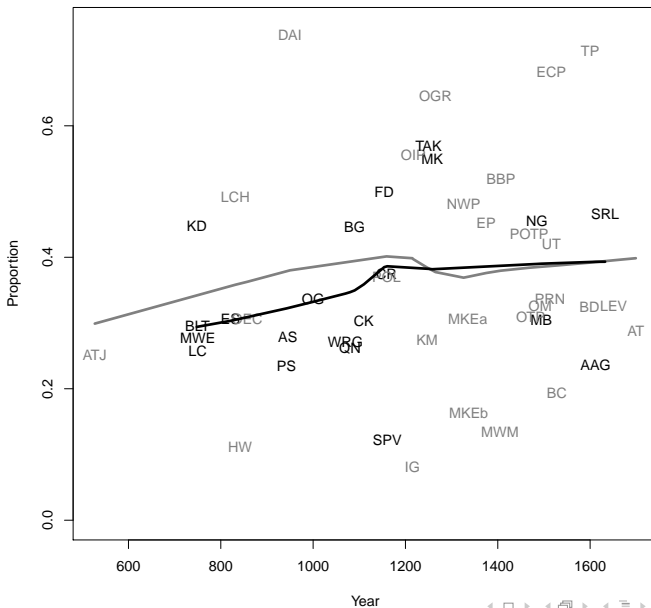
Practices and ideals of political rule

king, princ, citi, great, place, work, emperor, enemi, armi, letter

36.5% of paragraphs



# Coarse Topic 1



# Religion and Virtue

Connection between religion, virtue, justice and political rule



# Religion and Virtue

Connection between religion, virtue, justice and political rule

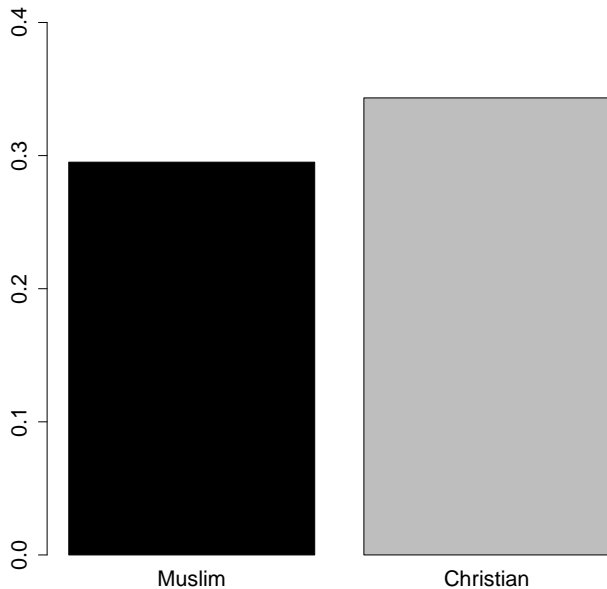
almighti,good,virtu,power,ruler,justic,prayer,rule,prophet,mena

# Religion and Virtue

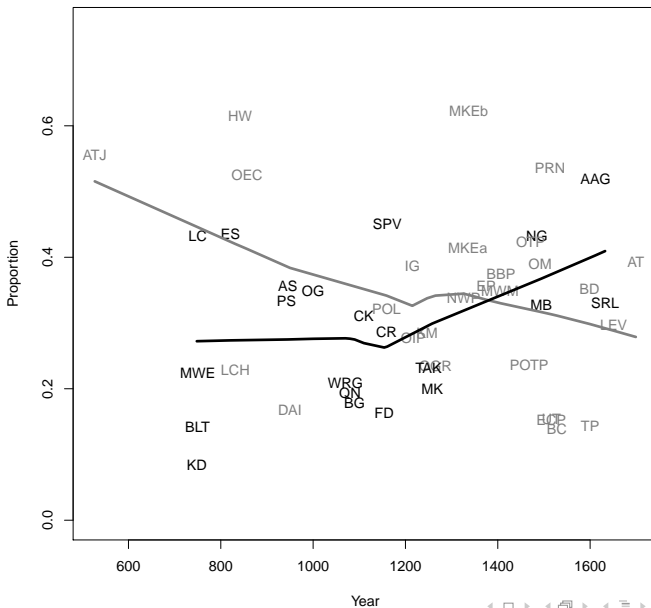
Connection between religion, virtue, justice and political rule

almighti,good,virtu,power,ruler,justic,prayer,rule,prophet,mena

32.2% of paragraphs



## Coarse Topic 2



# Inner Life of the Ruler

Personal relationships, care for and practices of the self, and ultimate fate of the soul

# Inner Life of the Ruler

Personal relationships, care for and practices of the self, and ultimate fate of the soul

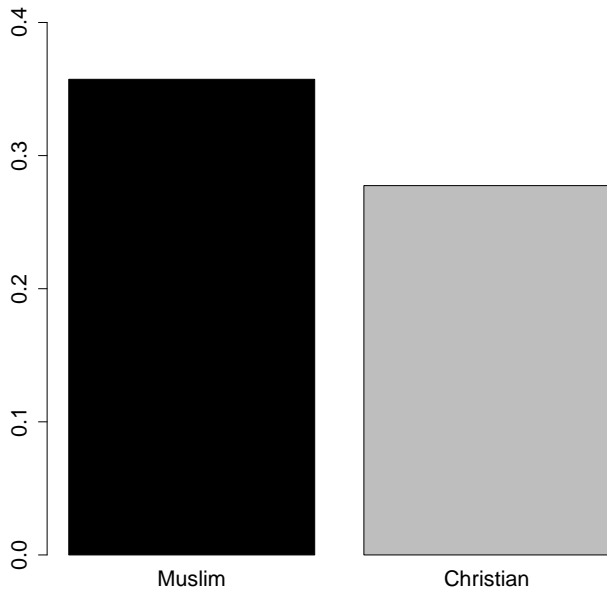
man,land,woman,know,bodi,eye,ladi,love,faculti,old

# Inner Life of the Ruler

Personal relationships, care for and practices of the self, and ultimate fate of the soul

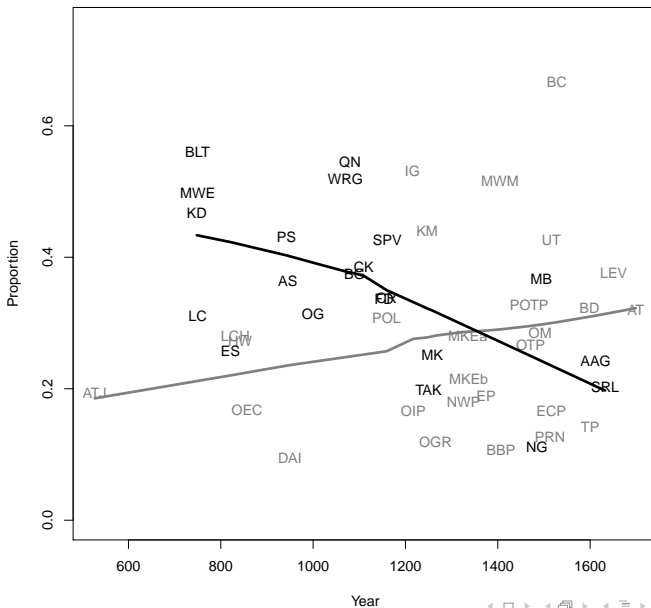
man,land,woman,know,bodi,eye,ladi,love,faculti,old

31.2% of paragraphs





# Coarse Topic 3



# Granular: Best Practices for Ruling

---

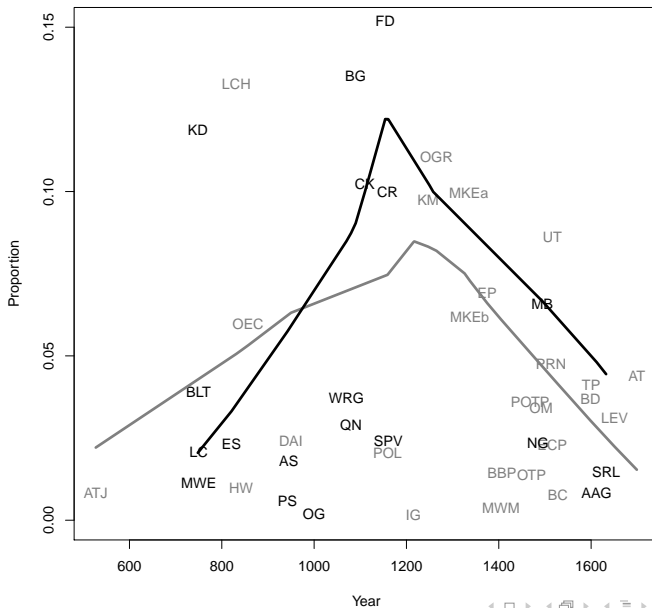
king, princ, citi, great, place, work, emperor, enemi, armi, letter

---

king, kingdom, royal, minist, reign, father, court, majesti, presenc, war

6.2% of paragraphs

# Coarse Topic 1 Granular Topic 1



# Granular: Characteristics that distinguish Just Ruler from Tyrant

---

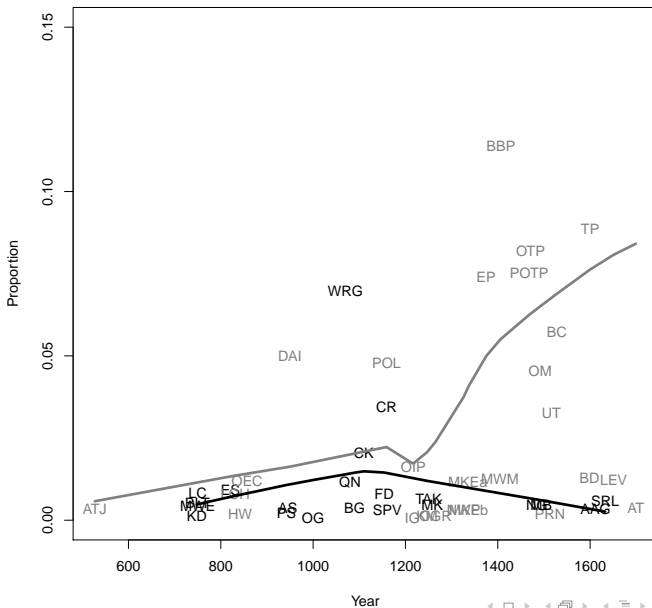
king, princ, citi, great, place, work, emperor, enemi, armi, letter

---

king, kingdom, royal, minist, reign, father, court, majesti, presenc, war  
princ, good, peopl, christian, tyranni, war, mind, ought, state, public

3.1% of paragraphs

## Coarse Topic 1 Granular Topic 2



# Granular: Religious Virtues and Political Ideals

---

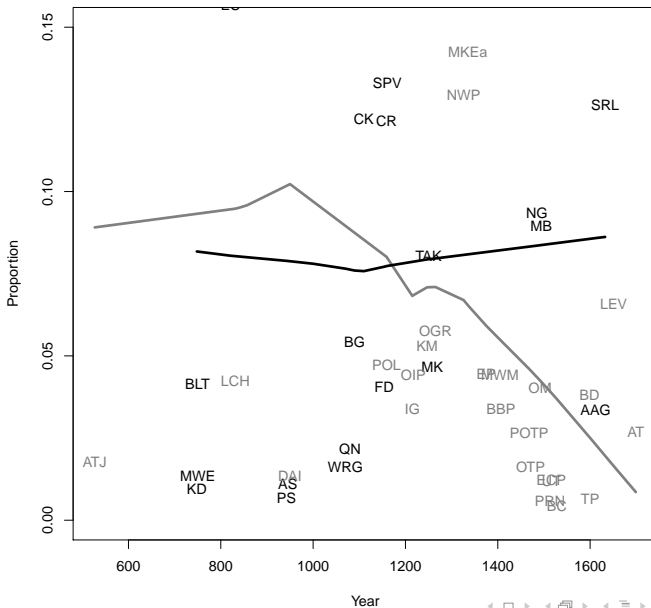
almighti,good,virtu,power,ruler,justic,prayer,rule,prophet,mena

---

almighti,bless,grant,peac,messeng,prophet,merci,holi,command,grace

6.9% of paragraphs

# Coarse Topic 2 Granular Topic 1



# Structural Topic Models

- Encode observed and unobserved meta data
- ?Improve substantive inferences

Next week:

- 1) Hanna Wallach on canonical topic models
- 2) Introduction to supervised learning

Work on your problem sets!



## Appendix: Inference for both Models

# Inference

Invariance in posterior makes it difficult (impossible) to approximate with sampling based methods (relabeling, aliasing problem).

Deterministic alternative: variational approximations.

Intuition: approximate posterior with simpler (still very general) approximating distribution.

Make approximation as “close” as possible

Approximate posterior with:

$$\begin{aligned} q(\alpha, \beta, \theta, \sigma, \pi, \tau) &= q(\alpha)q(\beta)q(\theta)q(\sigma)q(\pi)q(\tau) \\ &= q(\beta) \prod_{k=1}^K q(\theta)_k \prod_{i=1}^n \prod_{t=2005}^{2007} \left[ q(\sigma)_{it} q(\pi)_{it} \prod_{j=1}^J q(\tau)_{ijt} \right] \prod_{s=1}^S q(\alpha)_s \end{aligned}$$

# Variational Approximation

Optimization goal:

- Minimize the Kullback-Leibler divergence between approximating distribution  $q$  and true posterior  $p$ 
  - KL-divergence is a functional: takes **functions** as an input, returns a positive scalar
  - Measures “divergence” between two measures
- Use calculus of variations and theory from exponential models to derive iterative algorithm
- See “An Introduce to Bayesian Inference via Variational Approximations” for extended introduction (Grimmer, 2011)

# Minimize KL-divergence by Maximizing a Lower Bound

$$\log p(\mathbf{Y}) = \log \sum_{\sigma} \sum_{\tau} \iiint p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{Y}) d\theta d\alpha d\beta d\pi$$

$$\log p(\mathbf{Y}) = \log \sum_{\sigma} \sum_{\tau} \iiint p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{Y}) \frac{q(\alpha, \beta, \theta, \sigma, \pi, \tau)}{q(\alpha, \beta, \theta, \sigma, \pi, \tau)} d\theta d\alpha d\beta d\pi$$

$$\log p(\mathbf{Y}) \geq \underbrace{\sum_{\sigma} \sum_{\tau} \iiint q(\alpha, \beta, \theta, \sigma, \pi, \tau) \log \frac{p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{Y})}{q(\alpha, \beta, \theta, \sigma, \pi, \tau)} d\theta d\alpha d\beta d\pi}_{\mathcal{L}(q)}$$

# Minimize KL-divergence by Maximizing a Lower Bound

$$\log p(\mathbf{Y}) = \mathcal{L}(q) + \text{KL}(q||p)$$

# Minimize KL-divergence by Maximizing a Lower Bound

$$\begin{aligned}\log p(\mathbf{Y}) &= \mathcal{L}(q) + \text{KL}(q||p) \\ \underbrace{\log p(\mathbf{Y})}_{\text{fixed number}} &= \mathcal{L}(q) + \underbrace{\text{KL}(q||p)}_{\text{Positive}}\end{aligned}$$

# Minimize KL-divergence by Maximizing a Lower Bound

$$\begin{aligned}\log p(\mathbf{Y}) &= \mathcal{L}(q) + \text{KL}(q||p) \\ \underbrace{\log p(\mathbf{Y})}_{\text{fixed number}} &= \mathcal{L}(q) + \underbrace{\text{KL}(q||p)}_{\text{Positive}}\end{aligned}$$

If  $\mathcal{L}(q)$  get bigger,  $\text{KL}(q||p)$  get smaller.  $\Rightarrow$



# Minimize KL-divergence by Maximizing a Lower Bound

$$\begin{aligned}\log p(\mathbf{Y}) &= \mathcal{L}(q) + \text{KL}(q||p) \\ \underbrace{\log p(\mathbf{Y})}_{\text{fixed number}} &= \mathcal{L}(q) + \underbrace{\text{KL}(q||p)}_{\text{Positive}}\end{aligned}$$

If  $\mathcal{L}(q)$  get bigger,  $\text{KL}(q||p)$  get smaller.  $\Rightarrow$

If  $\mathcal{L}(q)$  is at a maximum  $\text{KL}(q||p)$  is at a minimum (duals).

# Maximizing $\mathcal{L}(q)$

Goal: choose  $q$  to maximize  $\mathcal{L}(q)$ .

$$q(\alpha)^{\text{old}}, q(\beta)^{\text{old}}, q(\theta)^{\text{old}}, q(\sigma)^{\text{old}}, q(\pi)^{\text{old}}, q(\tau)^{\text{old}}.$$

## Iterative Algorithm:

Choose,  $q(\pi)^{\text{new}}$  to max  $\mathcal{L}(q)$ —holding  $q(\theta)^{\text{old}}, q(\tau)^{\text{old}}, q(\alpha)^{\text{old}}, q(\beta)^{\text{old}}, q(\sigma)^{\text{old}}$  constant.

Choose,  $q(\theta)^{\text{new}}$  to max  $\mathcal{L}(q)$ —holding  $q(\pi)^{\text{new}}, q(\tau)^{\text{old}}, q(\alpha)^{\text{old}}, q(\beta)^{\text{old}}, q(\sigma)^{\text{old}}$  constant.

Choose,  $q(\tau)^{\text{new}}$  to max  $\mathcal{L}(q)$ —holding  $q(\theta)^{\text{new}}, q(\pi)^{\text{new}}, q(\alpha)^{\text{old}}, q(\beta)^{\text{old}}, q(\sigma)^{\text{old}}$  constant.

Choose,  $q(\alpha)^{\text{new}}$  to max  $\mathcal{L}(q)$ —holding  $q(\theta)^{\text{new}}, q(\tau)^{\text{new}}, q(\pi)^{\text{new}}, q(\beta)^{\text{old}}, q(\sigma)^{\text{old}}$  constant.

Choose,  $q(\beta)^{\text{new}}$  to max  $\mathcal{L}(q)$ —holding  $q(\theta)^{\text{new}}, q(\tau)^{\text{new}}, q(\pi)^{\text{new}}, q(\alpha)^{\text{new}}, q(\sigma)^{\text{old}}$  constant.

Choose,  $q(\sigma)^{\text{new}}$  to max  $\mathcal{L}(q)$ —holding  $q(\theta)^{\text{new}}, q(\tau)^{\text{new}}, q(\pi)^{\text{new}}, q(\alpha)^{\text{new}}, q(\beta)^{\text{new}}$  constant.

## Example for $q(\boldsymbol{\pi})^{\text{new}}$

$$\begin{aligned}\mathcal{L}(q) &= \int q(\boldsymbol{\pi})^{\text{new}} \left\{ \underbrace{\sum_{\boldsymbol{\sigma}} \sum_{\boldsymbol{\tau}} \iiint \log p(\mathbf{Y}, \boldsymbol{\pi}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \boldsymbol{\tau}) q(\boldsymbol{\sigma})^{\text{old}} q(\boldsymbol{\tau})^{\text{old}} q(\boldsymbol{\theta})^{\text{old}} q(\boldsymbol{\alpha})^{\text{old}} q(\boldsymbol{\beta})^{\text{old}} d\boldsymbol{\theta} d\boldsymbol{\alpha} d\boldsymbol{\beta}}_{\mathbb{E}_{\boldsymbol{\alpha}, \boldsymbol{\tau}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\beta}} [\log p(\mathbf{Y}, \boldsymbol{\pi}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \boldsymbol{\tau}, \boldsymbol{\beta}, \boldsymbol{\sigma})]} \right\} d\boldsymbol{\pi} \\ &\quad - q(\boldsymbol{\pi})^{\text{new}} \log q(\boldsymbol{\pi})^{\text{new}} + \text{constants}\end{aligned}$$

Define

$$\log \tilde{p}(\boldsymbol{\pi}) = \mathbb{E}_{\boldsymbol{\alpha}, \boldsymbol{\tau}, \boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\beta}} [\log p(\mathbf{Y}, \boldsymbol{\pi}, \boldsymbol{\alpha}, \boldsymbol{\theta}, \boldsymbol{\tau}, \boldsymbol{\beta}, \boldsymbol{\sigma})] + \text{constants}$$

## Example for $q(\pi)^{\text{new}}$

Substituting  $\log \tilde{p}(\pi)$ ,

$$\begin{aligned} &= \int q(\pi)^{\text{new}} \log \left( \frac{\tilde{p}(\pi)}{q(\pi)^{\text{new}}} \right) d\pi \\ &= -\text{KL}(q(\pi)^{\text{new}} \parallel \tilde{p}(\pi)) \end{aligned}$$

$\Rightarrow$  At a maximum when  $q(\pi)^{\text{new}} = \tilde{p}(\pi)$ .

Equivalently,

$$\begin{aligned} \log q(\pi)^{\text{new}} &= \log \tilde{p}(\pi) \\ &= E_{\alpha, \tau, \theta, \sigma, \beta} [\log p(\mathbf{Y}, \pi, \alpha, \theta, \tau, \beta, \sigma)] + \text{constants} \end{aligned}$$

Or,

$$q(\pi)^{\text{new}} = \frac{\exp(E_{\alpha, \tau, \theta, \sigma, \beta} [\log p(\mathbf{Y}, \pi, \alpha, \theta, \tau, \beta, \sigma)])}{\int \exp(E_{\alpha, \tau, \theta, \sigma, \beta} [\log p(\mathbf{Y}, \pi, \alpha, \theta, \tau, \beta, \sigma)]) d\pi}$$

To maximize  $\mathcal{L}(q)$  we use the following iterative algorithm

$$q(\sigma)^{\text{new}} \propto \exp(E_{\tau, \theta, \alpha, \beta, \pi}[\log p(\alpha, \beta, \theta, \sigma, \pi, \tau, Y)])$$

$$q(\tau)^{\text{new}} \propto \exp(E_{\sigma, \theta, \alpha, \beta, \pi}[\log p(\alpha, \beta, \theta, \sigma, \pi, \tau, Y)])$$

$$q(\theta)^{\text{new}} \propto \exp(E_{\sigma, \tau, \alpha, \beta, \pi}[\log p(\alpha, \beta, \theta, \sigma, \pi, \tau, Y)])$$

$$q(\alpha)^{\text{new}} \propto \exp(E_{\sigma, \tau, \theta, \beta, \pi}[\log p(\alpha, \beta, \theta, \sigma, \pi, \tau, Y)])$$

$$q(\beta)^{\text{new}} \propto \exp(E_{\sigma, \tau, \theta, \alpha, \pi}[\log p(\alpha, \beta, \theta, \sigma, \pi, \tau, Y)])$$

$$q(\pi)^{\text{new}} \propto \exp(E_{\sigma, \tau, \theta, \alpha, \beta}[\log p(\alpha, \beta, \theta, \sigma, \pi, \tau, Y)])$$

# Update for $q(\boldsymbol{\sigma})_{it}$

$q(\boldsymbol{\sigma})_{it}$  is a Multinomial( $1, \mathbf{c}_{it}$ ) distribution, with typical parameter  $c_{its}$

$$\mathbf{c}_{its} \propto \exp \left\{ \mathbb{E}[\log \beta_s] + \log \Gamma\left(\sum_{k=1}^K \alpha_{ks}\right) - \sum_{k=1}^K \log \Gamma(\alpha_{ks}) + \sum_{k=1}^K (\alpha_{ks} - 1) \mathbb{E}[\log \pi_{itk}] \right\}.$$

# Update for $q(\boldsymbol{\tau})_{ijt}$

$q(\boldsymbol{\tau})_{ijt}$  is a Multinomial( $1, \mathbf{r}_{ijt}$ ) distribution with typical parameter,

$$r_{ijtk} \propto \exp \left\{ E[\log \pi_{itk}] + \sum_{w=1}^W y_{ijtw} E[\log \theta_{kw}] \right\}.$$

# Update for $q(\boldsymbol{\pi})_{it}$

$q(\boldsymbol{\pi})_{it}$  is a Dirichlet( $\boldsymbol{\gamma}_{it}$ ) distribution, with typical parameter  $\gamma_{itk}$  equal to

$$\gamma_{itk} = \sum_{s=1}^S c_{its} \alpha_{sk}^* + \sum_{j=1}^{D_{it}} r_{ijtk}$$



# Update for $q(\boldsymbol{\theta})_k$

$q(\boldsymbol{\theta})_k$  is a Dirichlet( $\boldsymbol{\eta}_k$ ) distribution, with typical parameter equal to,

$$\eta_{kw} = \lambda_w + \sum_{i=1}^n \sum_{t=2005}^{2007} \sum_{j=1}^{D_{it}} r_{itjk} y_{itw}$$

# Update for $q(\beta)$

$q(\beta)$  is a Dirichlet( $\phi$ ) distribution, with typical parameter  $\phi_s$  equal to,

$$\phi_s = 1 + \sum_{i=1}^n \sum_{t=2005}^{2007} c_{its}$$

# Completing $q(\boldsymbol{\sigma})_{it}$ and $q(\boldsymbol{\tau})_{ijt}$

Finishing  $q(\boldsymbol{\sigma})_{it}$ :

- $E[\log \beta_s] = \Psi(\phi_s) - \Psi(\sum_{z=1}^S \phi_z)$  where  $\Psi(\cdot)$  is the digamma function (the derivative of the gamma function)
- $E[\log \pi_{itk}] = \Psi(\gamma_{itk}) - \Psi(\sum_{z=1}^K \gamma_{itz})$

Finishing  $q(\boldsymbol{\tau})_{ijt}$

- $E[\log \theta_{kw}] = \Psi(\eta_{kw}) - \Psi(\sum_{z=1}^W \eta_{kz})$ .

# Update Steps for $\alpha_s$

(Newton-Raphson, Minka 2000)

- Define  $N_s = \sum_{i=1}^n \sum_{t=2005}^{2007} c_{its}$ .

Differentiating with respect to  $\alpha_{ks}$  shows that

$$\frac{\partial \log q(\alpha)_k^{\text{new}}}{\partial \alpha_{ks}} = -\frac{1}{\lambda} + N_s \Psi\left(\sum_{k=1}^K \alpha_{ks}\right) - N_s \Psi(\alpha_{ks}) + \sum_{i=1}^n \sum_{t=2005}^{2007} c_{its} \frac{\left(\Psi(\gamma_{itk}) - \Psi\left(\sum_{z=1}^K \gamma_{itz}\right)\right)}{N_s}$$

- Call Gradient  $\frac{\partial \log q(\alpha)_k^{\text{new}}}{\partial \alpha_k}$ .
- Define H as the Hessian (matrix of second derivatives).
- Diagonal element  $h_{jj} = N_s \Psi'\left(\sum_{k=1}^K \alpha_{ks}\right) - N_s \Psi'(\alpha_{js})$  where  $\Psi'(\cdot)$  is the trigamma function
- Off-diagonal element ( $a \neq b$ )  $h_{ab} = N_z \Psi'\left(\sum_{k=1}^K \alpha_{ks}\right)$ .

For each  $s$  we iterate,

$$\alpha_s^{\text{new}} = \alpha_s^{\text{old}} - H^{-1} \frac{\partial \log q(\alpha)_k^{\text{new}}}{\partial \alpha_k}$$

until convergence

Initialize  $\gamma_{it}^{\text{old}}$  (for all  $i$  and  $t$ ),  $\eta_k^{\text{old}}$  (for all  $k$ ),  $\phi^{\text{old}}$ ,  $\alpha_s^{\text{old}}$  (for all  $s$ ).

Do until convergence in lower-bound.

- for all  $i, t, j$  and  $k$ , set

$$r_{ijk}^{\text{new}} \propto \exp \left( \Psi(\gamma_{itk}^{\text{old}}) - \Psi(\sum_{z=1}^K \gamma_{itz}^{\text{old}}) + \sum_{w=1}^W y_{ijtw} \left[ \Psi(\eta_{kw}^{\text{old}}) - \Psi(\sum_{z=1}^W \eta_{kz}^{\text{old}}) \right] \right)$$

-for all  $i, t$ , and  $s$  set

$$c_{its}^{\text{new}} \propto \exp(\Psi(\phi_s^{\text{old}}) - \Psi(\sum_{z=1}^S \phi_z^{\text{old}}) + \log \Gamma(\sum_{k=1}^K \alpha_{ks}^{\text{old}}) - \sum_{k=1}^K \log \Gamma(\alpha_{ks}^{\text{old}}) + \sum_{k=1}^K (\alpha_{ks}^{\text{old}} - 1) [\Psi(\gamma_{itk}^{\text{old}}) - \Psi(\sum_{z=1}^K \gamma_{itz}^{\text{old}})])$$

- for all  $i, t$ , and  $k$  set

$$\gamma_{itk}^{\text{new}} = \sum_{s=1}^S c_{its}^{\text{new}} \alpha_{ks}^{\text{old}} + \sum_{j=1}^{D_{it}} r_{ijk}^{\text{new}}$$

- for all  $k$  and  $w$  set

$$\eta_{kw}^{\text{new}} = \lambda_w + \sum_{i=1}^n \sum_{t=2005}^{2007} \sum_{j=1}^{D_{it}} r_{ijk}^{\text{new}} y_{ijtw}$$

-for all  $s$  set

$$\phi_s^{\text{new}} = 1 + \sum_{i=1}^n \sum_{t=2005}^{2007} c_{its}^{\text{new}}$$

- For all  $s$  obtain  $\alpha_s^{\text{new}}$  using Newton-Raphson algorithm.

- Evaluate lower-bound.

If converged:

Return posterior approximation.

## < / Variational Approximation >

## < Model Selection >

# Number of Topics

- 1) Substantive search (about 40-50)
- 2) 10-fold cross-validation. Loss function, approximate predictive posterior

$$p(\hat{\mathbf{y}}|\mathbf{Y}) \approx \sum_{\hat{\tau}} \iint p(\hat{\mathbf{y}}|\hat{\tau}, \boldsymbol{\theta}) p(\hat{\tau}|\boldsymbol{\pi}) q(\boldsymbol{\theta}|\boldsymbol{\eta}) q(\boldsymbol{\pi}) d\boldsymbol{\theta} d\boldsymbol{\pi}$$

- 3) Convergence with Nonparametric Bayesian model (Dirichlet process prior)

All converge on about 44 topics



# Number of Styles

$$\text{BIC} = 2 \log p(\mathbf{Y})$$

$$1) \text{ BIC} \approx 2(\mathcal{L}(q) + \log K! + \log S!) - (K \times S)(n)$$

$$2) \text{ BIC} \approx 2 \log p(\mathbf{Y} | \bar{\pi}, \bar{\theta}, \bar{\tau}) - (K \times S)(n)$$

## < / Model Selection >

## Hierarchy of Topics

# Posterior Distribution

$$\begin{aligned}
 p(\alpha, \pi, \eta, \beta, \sigma, \mu, \tau | \mathbf{X}) &\propto \prod_{m=1}^M c(\kappa) \exp\left(\kappa \boldsymbol{\eta}'_m \frac{\mathbf{1}}{\sqrt{J}}\right) \times \prod_{m=1}^M \prod_{k=1}^K \left[ \beta_m c(\kappa) \exp\left(\kappa \boldsymbol{\mu}'_k \boldsymbol{\eta}_m\right) \right]^{\sigma_{mk}} \prod_{k=1}^K \exp(-\alpha_k) \times \\
 &\prod_{i=1}^{48} \left[ \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K (\alpha_k)} \prod_{k=1}^K \pi_{ik}^{\alpha_k - 1} \times \prod_{j=1}^{D_i} \prod_{k=1}^K \left[ \pi_{ik} c(\kappa) \exp(\kappa \mathbf{x}_{ij}^* \boldsymbol{\mu}_k) \right]^{\tau_{ijk}} \right] \quad (0.1)
 \end{aligned}$$

Which we approximate with:

$$\begin{aligned}
 q(\alpha, \pi, \eta, \beta, \sigma, \mu, \tau) &= q(\alpha) q(\pi) q(\eta) q(\beta) q(\sigma) q(\mu) q(\tau) \quad (0.2) \\
 &= q(\alpha) \prod_{i=1}^{48} q(\pi)_i \prod_{m=1}^M q(\eta)_m q(\beta) \prod_{k=1}^K q(\sigma)_k \prod_{k=1}^K q(\mu)_k \prod_{i=1}^{48} \prod_{j=1}^{D_i} q(\tau)_{ij}
 \end{aligned}$$

## Update for $q(\boldsymbol{\sigma})_k$

$q(\boldsymbol{\sigma})_k$  is a Multinomial( $1, \mathbf{c}_k$ ) where typical element  $c_{mk}$  is equal to

$$c_{mk} \propto \exp(E[\log \beta_m] + E[\kappa \boldsymbol{\mu}_k \boldsymbol{\eta}_m]) .$$

We will complete the update step when we have determined the remaining forms of the distribution

# Update for $q(\boldsymbol{\tau})_{ij}$

$q(\boldsymbol{\tau})_{ij}$  is a Multinomial(1,  $\mathbf{r}_{ij}$ , with typical element of  $r_{ijk}$  equal to

$$r_{ijk} \propto \exp \left( \mathbb{E}[\log \pi_{ik}] + \mathbb{E}[\kappa \mathbf{y}_{ij}^* \boldsymbol{\mu}_k] \right)$$

Again, as we complete the parametric forms of the other update steps we can complete this update equation.

## Update for $q(\boldsymbol{\pi})_i$

$q(\boldsymbol{\pi})_i$  is a Dirichlet( $\boldsymbol{\gamma}_i$ ) distribution, where typical element  $\gamma_{ik}$  is equal to

$$\gamma_{ik} = \alpha_k + \sum_{j=1}^{D_i} r_{ijk}$$

# Update for $q(\beta)$

$q(\beta)$  is a Dirichlet( $\phi$ ) distribution with typical parameter  $\phi_m$  equal to

$$\phi_m = 1 + \sum_{k=1}^K c_{mk}$$



## Update for $q(\boldsymbol{\eta})_m$

Given the complications of taking expectations with the vMF distribution, we instead provide maximization steps for the vMF parameters. To obtain the form of the updates we follow the derivation outlined in Banerjee et al (2005). To do this, we take the log of the posterior distribution and identify the parameters that depend upon  $\boldsymbol{\eta}_m$ .

$$\log(p(\boldsymbol{\eta}_m)) = \sum_{k=1}^K c_{km} \kappa \boldsymbol{\mu}_k \boldsymbol{\eta}_m + \kappa \boldsymbol{\eta}_m \frac{1}{\sqrt{J}} + \text{constants}$$

## Update for $q(\boldsymbol{\eta})_m$

To set up the constrained optimization we also introduce the Lagrangian  $\lambda$ , with the constraint that  $\boldsymbol{\eta}_m' \boldsymbol{\eta}_m = 1$ ,

$$\log(p(\boldsymbol{\eta}_m)) \propto \sum_{k=1}^K c_{km} \kappa \boldsymbol{\mu}_k \boldsymbol{\eta}_m + \kappa \boldsymbol{\eta}_m \frac{1}{\sqrt{J}} - \lambda(\boldsymbol{\eta}_m' \boldsymbol{\eta}_m - 1).$$

Differentiating with respect to  $\boldsymbol{\eta}_m$ , setting equal to zero and solving yields

$$\frac{\kappa \left( \sum_{k=1}^K c_{mk} \boldsymbol{\mu}_k + \frac{1}{\sqrt{J}} \right)}{2\lambda} = \boldsymbol{\eta}_m \quad (0.3)$$

## Update for $q(\boldsymbol{\eta})_m$

If we differentiate with respect to  $\lambda$  and solve we see that  $\boldsymbol{\eta}'_m \boldsymbol{\eta}_m = 1$  or that  $\|\boldsymbol{\eta}'_m \boldsymbol{\eta}_m\| = 1$ . Substituting this into Equation 0.3 we have,

$$\begin{aligned} \frac{\kappa}{2\lambda} \left( \left( \sum_{k=1}^K c_{mk} \boldsymbol{\mu}_k + \frac{\mathbf{1}}{\sqrt{J}} \right)' \left( \sum_{k=1}^K c_{mk} \boldsymbol{\mu}_k + \frac{\mathbf{1}}{\sqrt{J}} \right) \right)^{1/2} &= 1 \\ \frac{\kappa \left\| \sum_{k=1}^K c_{mk} \boldsymbol{\mu}_k + \frac{\mathbf{1}}{\sqrt{J}} \right\|}{2} &= \lambda \end{aligned}$$

Doing a final substitution we have

$$\boldsymbol{\eta}_m^* = \frac{\sum_{k=1}^K c_{mk} \boldsymbol{\mu}_k + \frac{\mathbf{1}}{\sqrt{J}}}{\left\| \sum_{k=1}^K c_{mk} \boldsymbol{\mu}_k + \frac{\mathbf{1}}{\sqrt{J}} \right\|}$$

# Update for $q(\boldsymbol{\mu})_k$

Following a very similar set of derivations, the update step for  $\boldsymbol{\mu}_k$  is

$$\boldsymbol{\mu}_k^* = \frac{\sum_{i=1}^{48} \sum_{j=1}^{D_i} r_{ijk} \mathbf{x}_{ij}^* + \sum_{m=1}^M c_{mk} \boldsymbol{\eta}_m^*}{\|\sum_{i=1}^{48} \sum_{j=1}^{D_i} r_{ijk} \mathbf{x}_{ij}^* + \sum_{m=1}^M c_{mk} \boldsymbol{\eta}_m^*\|}$$

# Completing updates for $q(\boldsymbol{\sigma})_k$ and $q(\boldsymbol{\tau})_{ij}$

Given the forms  $E[\log \beta_m] = \Psi(\phi_m) - \Psi(\sum_{m=1}^M \phi_m)$  and  $E[\log \pi_{ik}] = \Psi(\gamma_{ik}) - \Psi(\sum_{k=1}^K \gamma_{ik})$  where  $\Psi(\cdot)$  is the Digamma function.

# Update for $q(\alpha)$

A closed form update for the  $\alpha$  parameters is unavailable. So we use the Newton-Raphson algorithm outlined in Minka (2000) and Blei, Ng, and Jordan (2003).